# Multivariate error covariance estimates by Monte-Carlo simulation for assimilation studies in the Pacific Ocean.

Anna Borovikov*
*SAIC, Beltsville, Maryland*

Michele M. Rienecker
*Global Modeling and Assimilation Office,*
*NASA/Goddard Space Flight Center,*
*Greenbelt, Maryland*

Christian L. Keppenne
*SAIC, Beltsville, Maryland*

Gregory C. Johnson
*NOAA/Pacific Marine Environmental Laboratory*
*Seattle, Washington*

June 30, 2004

*\*Corresponding author address:* Anna Borovikov, Code 900.3 NASA/Goddard, Greenbelt, MD 20771, *ayb@mohawk.gsfc.nasa.gov*

**Abstract**

One of the most difficult aspects of ocean state estimation is the prescription of the model forecast error covariances. The paucity of ocean observations limits our ability to estimate the covariance structures from model-observation differences. In most practical applications, simple covariances are usually prescribed. Rarely are cross-covariances between different model variables used. Here a comparison is made between a univariate Optimal Interpolation (UOI) scheme and a multivariate OI algorithm (MvOI) in the assimilation of ocean temperature profiles. In the UOI case only temperature is updated using a Gaussian covariance function. In the MvOI, salinity, zonal and meridional velocities as well as temperature are updated using an empirically estimated multivariate covariance matrix.

Earlier studies have shown that a univariate OI has a detrimental effect on the salinity and velocity fields of the model. Apparently, in a sequential framework it is important to analyze temperature and salinity together. For the MvOI an estimation of the model error statistics is made by Monte-Carlo techniques from an ensemble of model integrations. An important advantage of using an ensemble of ocean states is that it provides a natural way to estimate cross-covariances between the fields of different physical variables constituting the model state vector, at the same time incorporating the model's dynamical and thermodynamical constraints as well as the effects of physical boundaries.

Only temperature observations from the Tropical Atmosphere-Ocean array have been assimilated in this study. In order to investigate the efficacy of the multivariate scheme, two data assimilation experiments are validated with a large independent set of recently published subsurface observations of salinity, zonal velocity and temperature. For reference, a control run with no data assimilation is used to check how the data assimilation affects systematic model errors.While the performance of the UOI and MvOI is similar with respect to the temperature field, the salinity and velocity fields

are greatly improved when the multivariate correction is used, as evident from the analyses of the rms differences between these fields and independent obsevations. The MvOI assimilation is found to improve upon the control run in generating water masses with properties close to the observed, while the UOI failed to maintain the temperature and salinity structure.

## 1.  Introduction

Data assimilation provides a framework for the combination of the information about the state of the ocean contained in an incomplete data stream with our knowledge of the ocean dynamics included in a model. The problem of data assimilation may be formulated in statistical terms where, because of uncertainty in both observations and models, an estimate of the state of the ocean at any given time is considered to be a realization of a random variable.    An estimate of the state of the ocean is produced as a blend of observation and model estimates based on prior knowledge of the error statistics of each, with some measure of the uncertainty in the estimate. The differences between assimilation methods lie primarily in the approaches taken to estimate the error statistics associated with the forward (dynamical) model, the so-called background or forecast error statistics. Since an accurate representation of the observation and model error statistics is crucial to a successful data assimilation, a lot of effort has been expended in this direction.

One simplifying assumption that is often made is that the forecast error statistics do not change significantly with time and thus can be approximated by a constant probability distribution. This is the basis of the Optimal Interpolation (OI) data assimilation scheme, also known as statistical interpolation (e.g., Daley 1991, chapters 4 and 5). An alternative to this assumption is to allow for time evolution of the probability distribution. An example of such a data assimilation scheme is the Kalman Filter (Kalman 1960), in which the model and data errors are assumed to be normally distributed and the forecast error covariance matrix is evolved prognostically. The Kalman Filter can be shown to give an optimal estimate in the case of linear dynamics and linear observation operator. To account for nonlinear processes a generalization of the Kalman Filter, the Extended Kalman Filter uses instantaneous linearization (and often a truncation) of the model equations during the update of the error covariance matrix and the full equations to update the model (e.g., Daley 1991; Ghil and Malanotte-Rizzoli 1991). However, time stepping the forecast error covariance matrix is computationally expensive, rendering this method impractical when used

with high-resolution general circulation models. Under certain conditions it is possible to use an asymptotic Kalman Filter (e.g., Fukumori et al. 1993), where a steady-state covariance matrix replaces the time-evolving one. An Ensemble Kalman Filter (EnKF) was introduced by Evensen (1994) based on a Monte Carlo technique, in which the forecast error statistics are computed from an ensemble of model states evolving simultaneously. The methodology of the EnKF was further refined by adding pertubations to the observations (e.g., Burgers et al. 1998) to maintain consistent variance in the ensemble analysis. An application of this method with the Poseidon ocean model used in this study has been developed by Keppenne and Rienecker (2002, 2003). Zhang and Anderson (2003) describe an adjustment Kalman filter (EAKF) which is another modification of the Kalman filter based on a Monte Carlo approach, and compare it to an ensemble OI scheme (time-invariant forecast error, but spatial structure is derived from a collection of state vectors) as well as an OI with functionally prescribed covariances. Their conclusion is that when applied to a simple atmospheric model an ensemble OI can produce reasonably good assimilation results if the covariance matrix is chosen appropriately.

This study focuses on the importance of the multivariate aspect of the forecast error covariance in the context of data assimilation using OI. Provided with a fairly good observing network, the background error structure can be estimated using analysis of spatial and temporal decorrelation scales, as done in numerous meteorological applications (Ghil and Malanotte-Rizzoli 1991; Derber et al. 1991). Several studies used a Monte Carlo approach to estimate forecast error covariance structure from an ensemble of assimilation integrations with perturbed models and randomly selected (Buehner, 2004) or randomly perturbed (Houtekamer, 1996) observations. However, even for atmospheric data assimilation, the observing system is not adequate to support a full calculation of background error covariance statistics, hence model forecasts are often used for error estimation, as, for example, done in NMC assimilation algorithm (Derber and Parrish 1991).

The vastness and complexity of the domain and relative scarcity of oceanographic obser-

vations would require additional simplifying assumptions in similar calculations. To avoid imposing severe restrictions on the error covariance calculation due to limited data availability, this paper explores the efficacy of estimating the forecast error from an ensemble of model integrations. A Monte-Carlo technique similar to the EnKF is used here. An important advantage of using an ensemble of ocean states is that it provides a natural way to estimate cross-covariances between the fields of different physical variables constituting the model state vector while incorporating model balance relations and the influence of boundaries. The idea of a multivariate forecast error covariance matrix has been implemented in the oceanographic context, for example, to relate the tide gauge data (Cane et al. 1996) and surface velocity data (Oke at al. 2002) to the dynamically varying quantities in the water column below.

There are many questions that arise with this approach. For example, how large should the ensemble be, and more generally, how should it be generated? Other questions are related to the underlying assumption of the stationarity and the unbiased nature of error statistics in the OI algorithm. Will a one-time estimate of the forecast error, derived from a Monte Carlo ensemble, be a good representation of this error at another time, at any time during assimilation? Or, in other words, what is the variability of the forecast error covariance structure? What are the dominant time scales? Can this information be acquired and, if so, used to improve the assimilation scheme?

The primary interest of this study is ocean phenomena taking place on seasonal-to-interannual time scales. One example of such phenomena is the quasi-regular occurrence of El Niño - a large scale warming of near-surface temperature in the eastern equatorial Pacific Ocean accompanied by a basin wide perturbation in the tilt of the thermocline across the equatorial ocean (e.g., Philander 1990). The estimate of error statistics derived below attempts to capture errors associated with such variability. The logical organization of the paper is as follows. Next the OI assimilation algorithm, model and data are described (Section 2). Then the forecast error covariance model, a traditional Gaussian model of

the forecast error covariance and the empirical multivariate model of interest are detailed (Section 3). Then the multivariate error covariance model properties are explored (Section 4). After the experimental setup is decribed, the results of multivariate assimilation are compared with univariate assimilation (Section 5). The paper concludes with discussion of the results and further directions of research (Section 6).

## 2. OI assimilation

### a. OI framework

A detailed discussion of the sequential data assimilation algorithms can be found in earlier literature (see for example, Lorenc (1988), Daley (1991) or Cohn (1997)). Here, only a brief outline is given to inroduce necessary terminology and notation.

The aim of a data assimilation algorithm is to determine the best estimate of the state vector based on the estimates available from both model and observations. A dynamic (prediction) model can be represented in terms of a nonlinear operator $\Psi(\mathbf{x})$, where $\mathbf{x}$ is a state vector of length $n_x$. Let $\mathbf{d}$ denote a vector of observations which has dimension $n_d \ll n_x$ (typically for ocean applications) and an element of $\mathbf{d}$ is not necessarily an element of the state vector $\mathbf{x}$. A discrete form of the model can be written as $\mathbf{x}_k = \Psi_{k-1}(\mathbf{x}_{k-1})$, where $\mathbf{x}_k$ is the forecast state vector at time level $k$ and $\Psi_{k-1}$ is the numerical approximation to the set of model equations describing the evolution of the state forward from time $k-1$ to $k$. Similarly, observations available at time $k$ can be denoted as $\mathbf{d}_k$ and the observation transformation operator as $\mathcal{H}_k(\mathbf{x}_k)$.

A sequential, unbiased assimilation scheme for the time-varying $\mathbf{x}_k$ is given by:

$$\mathbf{x}_k^f = \Psi_{k-1}(\mathbf{x}_{k-1}^a) \tag{1}$$

$$\mathbf{x}_k^a = \mathbf{x}_k^f + \mathbf{K}_k \left( \mathbf{d}_k - \mathcal{H}_k(\mathbf{x}_k^f) \right) \tag{2}$$

Here superscript $f$ stands for the forecast and $a$ for the analysis. The sequential data assimilation schemes that have the form of equation (2) differ from each other by the weight

4

matrix $\mathbf{K}_k$ often called the *gain matrix*.

The optimality of $\mathbf{K}_k$ can be defined under certain assumptions about the error statistics. Most sequential data assimilation algorithms are based on assumptions that the observational and model errors are unbiased, white in time, spatially uncorrelated with each other and that their spatial covariances are known (usually it is assumed that, at least initially, the errors are Gaussian). The observational error may also include any error of representation of the processes of interest, although such errors will not in general satisfy the assumption of a white, Gaussian sequence. Without any loss of generality, it is also assumed that the system noise and the observational noise are uncorrelated with each other. Under these assumptions, for a linear model and a linear observation transformation operator, $\mathcal{H}_k \equiv \mathbf{H}_k$, the optimal $\mathbf{K}_k$ is given by

$$\mathbf{K}_k = \mathbf{P}_k^f \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^f \mathbf{H}_k^T + \mathbf{R}_k)^{-1}. \tag{3}$$

Here $\mathbf{P}_k^f$ is the forecast error covariance matrix, which, in general, is time-dependent. For a high resolution ocean model with the number of state variables on the order of $10^6$, $\mathbf{P}_k^f$ is extremely expensive to store and evaluate in full. Thus, numerous approaches have been suggested to simplify the computation of $\mathbf{P}_k^f$. The traditional OI method assumes that $\mathbf{P}_k^f \equiv \mathbf{P}$ is approximately constant in time. In the case of observational errors, the matrix $\mathbf{R}$ is often assumed to be diagonal and to contain only information about the level of variance in the measurement error due to instrumental imperfection and unresolved small-scale signal. There are means of allowing for simple time evolution of the forecast error variance (see, for example, Ghil and Malanotte-Rizzoli 1991; Rienecker and Miller 1991), but they are not considered here. A full evolution of $\mathbf{P}_k^f$ would be a Kalman filter.

The effects of non-linear dynamics and inhomogeneities associated with ocean boundaries are implicitly taken into account when the empirical forecast error covariance matrix $\mathbf{P}$ is constructed from model integrations as presented in the next section.

*b. Model and forcing*

The model used for this study is the Poseidon reduced-gravity, quasi-isopycnal ocean model introduced by Schopf and Loughe (1995) and used by Keppenne and Rienecker (2002, 2003) for testing the Ensemble Kalman Filter. The model described by Schopf and Loughe (1995) has been updated to include the effects of salinity (e.g., Yang et al. 1999). The model was shown to provide realistic simulations of tropical Pacific climatology and variability (Borovikov et al. 2001). Explicit detail of the model are provided in Schopf and Loughe (1995). The prognostic variables are layer thickness, temperature, salinity and the zonal and meridional current components. The generalized vertical coordinate of the model includes a turbulent well-mixed surface layer with entrainment parameterized according to a Kraus-Turner (1967) bulk mixed layer model.

For this study, the domain is restricted to the Pacific Ocean (45°S to 65°N) with realistic land boundaries. At the southern boundary the model temperature and salinity are relaxed to the Levitus (1994) climatology. The horizontal resolution of the model is 1° in longitude; and in the meridional direction a stretched grid is used, varying from 1/3° at the equator to 1° poleward of 10°S and 10°N. The calculation of the effects of vertical diffusion, implemented at three-hour intervals through an implicit scheme, are parameterized using a Richardson number-dependent vertical mixing following Pacanowski and Philander (1981). The diffusion coefficients are enhanced when needed to simulate convective overturning in cases of gravitationally unstable density profiles. Horizontal diffusion is also applied daily using an 8th-order Shapiro (1970) filter. The net surface heat flux is estimated using the atmospheric mixed layer model of Seager et al. (1994) with monthly averaged time-varying air temperature and specific humidity from the NCEP-NCAR reanalysis (e.g., Kalnay et al. 1996) and climatological shortwave radiation from the Earth Radiation Budget Experiment (ERBE) (e.g., Harrison et al. 1993), and climatological cloudiness from the International Satellite Cloud Climatology Project (ISCCP) (e.g., Rossow and Schiffer 1991).

Surface wind stress forcing is obtained from the Special Sensor Microwave Imager (SSM/I)

surface wind analysis (Atlas et al. 1991) based on the combination of the Defense Meteorolog-
ical Satellite Program (DMSP) SSM/I data with other conventional data and the ECMWF
10m surface wind analysis. The surface stress was produced from this analysis using the
drag coefficient of Large and Pond (1982). Monthly averaged wind stress forcing was applied
to the model. The precipitation is given by monthly averaged analyses of Xie and Arkin
(1997).

Model mean (1988-1997) temperature, salinity and zonal velocity sections along the equa-
tor compare very well with estimates made from observations (Johnson et al. 2002) taken
during an overlapping period (figure 1).

*c.   Data*

The TAO/Triton Array (figure 2), consisting of more than 70 moored buoys spanning
the equatorial Pacific (http://www.pmel.noaa.gov/toga-tao/home.html and McPhaden et
al. 1998), measures oceanographic and surface meteorological variables: air temperature,
relative humidity, surface winds, sea surface temperatures and subsurface temperatures down
to a depth of 500 meters. By 1994 these measurements became available daily approximately
uniformly spaced at 10-15° longitude and 2-3° latitude degrees across the equatorial Pacific
Ocean.

The temperature observations from the TAO/Triton array were the only data type used
during these assimilation experiments since the focus is on well-known deleterious effects of
temperature assimilation in the equatorial waveguide, as discussed, for example, in Troccoli
et al. 2002 and in Troccoli et al. 2003. [In the global assimilation conducted by the NASA
Seasonal-to-Interannual Prediction Project to initialize seasonal forecasts, the global XBT
data base is included.] The standard deviation of the observational error, denoted $\sigma_{TAO}$, is
set to 0.5°C and the errors are assumed to be uncorrelated in space and time. This value
is high compared to the instrumental error of 0.1°C (Freitag et al. 1994) since it also has
to reflect the representativeness error - i.e., the data contains a mixture of signals of various

7

scales including frequencies much higher than the target scales of assimilation. By tuning $\sigma_{TAO}$ we effectively control the ratio of the data error variance to the model error variance.

## 3. Forecast error covariance modeling

In error covariance structure modeling, one is striving for an accurate representation of the error statistics as well as for simple and efficient implementation for computational viability. With little knowledge of the true nature of the model error covariances, one often has to make assumptions and settle for simple methods that usually have the advantage of being easy to implement. This section describes two different models for the forecast error covariance structure, a simpler and less computationally intense and a more elaborate and more accurate model. For both, an OI framework is used wherein the forecast error covariance matrix, $\mathbf{P}^f$, is assumed to be time-invariant.

### a.  Univariate functional model

A commonly used analytical error covariance function (e.g., Carton and Hackert 1990 and Ji et al. 1995) has been employed here in the tropical Pacific Ocean region: the spatial structure of the model temperature (T) forecast error is assumed to be Gaussian in all three dimensions with scales 15°, 4° and 50 m in zonal, meridional and vertical directions, respectively. The values used in this study were estimated from the ensemble of model integrations described in the next subsection. Those spatial scales are also (marginally) resolved by the equatorial moorings which are nominally separated by 10° to 15° in the zonal direction and by 2° to 3° in the meridional direction. Horizontal scales are comparable to scales used in similar (3DVar) assimilation schemes (e.g., Ji et al. 1995 and Rosati et al. 1996). There are several advantages to this error covariance model. For the Gaussian form of the covariance function, the minimum variance estimate for the least squares minimizing functional is the maximum likelihood estimate, and the analysis error covariance function is also Gaussian.

8

It is relatively easy to implement and adapt to parallel computing architecture. The study by Rosati et al. (1997) also shows that use of such empirical covariance scales, though simplified, are nevertheless effective for improving seasonal forecasts.

In the present implementation the temperature observations have been processed and the correction was only made to the model temperature field during each assimilation cycle, while other variables adjusted according to the model's dynamic response to the temperature correction.

### b. Monte Carlo method for estimating the multivariate forecast error covariance

A more realistic covariance structure that is consistent with model dynamics and the presence of ocean boundaries was sought through an application of the Monte Carlo method. The variability across an ensemble of ocean state estimates was used for a one-time estimate of the model forecast error statistics. This approach is similar in spirit to the Ensemble Kalman Filter except that the error covariance does not evolve with time and does not feel the impact of prior data assimilation, although it could.

The design of this forecast error covariance model was influenced by the need to assimilate TAO mooring observations for seasonal forecasts. While the Poseidon model has a layered configuration, the TAO observations are taken at approximately constant depth levels. In the implementation for this study, the covariances are calculated on pre-defined depth levels. At each assimilation cycle the model fields are interpolated to these depths, the assimilation increments are computed on these pre-specified levels, and are then interpolated back to the temperature grid points at the center of the model layers. The discussion below deals with the three-dimensional model error covariance matrix whose horizontal structure coincides with the model grid, and in the vertical is arranged at depths coincident with the nominal TAO instrument depths.

9

Consider the non-dimensionalized model state vector

$$\mathbf{x} = \begin{bmatrix} T/\sigma_T \\ S/\sigma_S \\ U/\sigma_U \\ V/\sigma_V \\ ssh/\sigma_{ssh} \end{bmatrix}, \tag{4}$$

here $T$, $S$, $U$, $V$ and $ssh$ are model variables: temperature, salinity, zonal and meridional velocities and dynamic height respectively, and $\sigma_{[T,S,U,V,ssh]}$ are non-dimensionalizing factors. For the latter we took the global standard deviation within each of the model fields at a depth of 100 m (the depth of highest variability, around the thermocline): $\sigma_T{=}0.65$, $\sigma_S{=}0.08$, $\sigma_U{=}0.09$, $\sigma_V{=}0.08$ and $\sigma_{ssh}{=}0.08$ in the corresponding units. Note that $\sigma_T$, which represents the internal variability of the model is comparable to assumed $\sigma_{TAO}$ - the observational error standard deviation, so that the model and data are given comparable weights in assimilation. The multivariate covariance matrix is

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}^{T,T} & \mathbf{P}^{T,S} & \mathbf{P}^{T,U} & \mathbf{P}^{T,V} & \mathbf{P}^{T,ssh} \\ \mathbf{P}^{T,S} & \mathbf{P}^{S,S} & \mathbf{P}^{S,U} & \mathbf{P}^{S,V} & \mathbf{P}^{S,ssh} \\ \mathbf{P}^{U,T} & \mathbf{P}^{U,S} & \mathbf{P}^{U,U} & \mathbf{P}^{U,V} & \mathbf{P}^{U,ssh} \\ \mathbf{P}^{V,T} & \mathbf{P}^{V,S} & \mathbf{P}^{V,U} & \mathbf{P}^{V,V} & \mathbf{P}^{V,ssh} \\ \mathbf{P}^{ssh,T} & \mathbf{P}^{ssh,S} & \mathbf{P}^{ssh,U} & \mathbf{P}^{ssh,V} & \mathbf{P}^{ssh,ssh} \end{bmatrix}. \tag{5}$$

If the matrix $\mathbf{A}^{m \times n_x}$ contains the $m$-member ensemble of (anomalous) ocean states as columns, then $\mathbf{P}$ can be computed as

$$\mathbf{P}^{n_x \times n_x} = \frac{\mathbf{A}\mathbf{A}^T}{m-1}, \text{ with } rank(\mathbf{P}) \leq \min\{m, n_x\}. \tag{6}$$

The size of $\mathbf{P}$ is on the order of $n_x \approx 10^6$ (the dimension of the state vector), while its rank is smaller than the size of the ensemble, $m$ (on the order of $10^2$ in the case of this study). The estimated error covariance matrix was stored on file and read in during every assimilation cycle of the OI algorithm. Since the rank of the error covariance matrix $\mathbf{P}$ estimated using this method is no larger than the Monte Carlo ensemble size, it can be conveniently represented using a basis of empirical-orthogonal functions (eofs), $\mathbf{E}$. Eofs have been widely employed in oceanographic contexts, and a relevant theoretical background can be found, for example, in Preisendorfer (1988).

10

To compute the eof representation of $\mathbf{P}$, observe that $\mathbf{A}\mathbf{A}^T$ has the same eigenvalues as $\mathbf{A}^T\mathbf{A}$, which is only $m \times m$ and the eigenvectors of $\mathbf{A}\mathbf{A}^T$ are related to those of $\mathbf{A}^T\mathbf{A}$ as

$$\mathbf{E} = \mathbf{A}\mathbf{U}(\Lambda)^{-1/2}, \tag{7}$$

where $\mathbf{E}^{n_x \times m}$ contains the eigenvectors of $\mathbf{A}\mathbf{A}^T$, $\mathbf{U}^{m \times m}$ contains the eigenvectors of $\mathbf{A}^T\mathbf{A}$ and $\Lambda^{m \times m} = diag(\lambda_1^2, ..., \lambda_m^2)$ has the eigenvalues of $\mathbf{A}^T\mathbf{A}$. Then

$$\mathbf{P} = \frac{\mathbf{A}\mathbf{A}^T}{m-1} = \frac{\mathbf{E}\Lambda\mathbf{E}^T}{m-1} = \mathbf{L}\mathbf{L}^T. \tag{8}$$

The columns of $\mathbf{E}$ are orthonormal and the eigenvalues $\lambda_i^2$, $i = 1, .., m$, are the variances. Equation (3) can thus be rewritten as

$$\mathbf{K} = \mathbf{L}\mathbf{L}^T\mathbf{H}^T(\mathbf{H}\mathbf{L}\mathbf{L}^T\mathbf{H}^T + \mathbf{R})^{-1}, \text{ with } \mathbf{L} = \mathbf{E}\Lambda^{1/2}(m-1)^{-1/2}. \tag{9}$$

*1)   Ensemble generation*

As the first test of this methodology, the ensemble of states was generated by forcing the ocean model with an ensemble of air-sea fluxes:

$$\mathbf{F}_n = \mathbf{F} + \delta\mathbf{F}_n. \tag{10}$$

$\mathbf{F}$ is the forcing used for the control run, $\delta\mathbf{F}_n$ are interannual anomalies - in phase with respect to the annual cycle and interannual SST anomalies but with different internal atmospheric chaotic variations. Surface forcing is used for the ensemble generation because this is probably the dominant source of error in the upper ocean in the equatorial Pacific. Our approach is similar to Cane et al. (1996) in the sense that all the ensemble variability is a result of the perturbations to the atmospheric forcing, although the implementation details differ. Although errors in the synoptic forcing will be large, the focus here is on the longer time scales of interest for seasonal prediction. The fluxes were obtained from a series of integrations of the Aries atmospheric model (e.g., Suarez and Takacs 1995) forced by the same interannually varying sea surface temperatures (SST) and differing only in slight perturbations to the initial atmospheric state. The interannual anomalies in surface stress and

heat flux components were added to seasonal forcing estimated from the sources described in the section 2(b). This approach attributes all of the ocean model forecast error to uncertainties in the surface flux anomalies, since differences between the ensemble members were due to atmospheric internal variability. No perturbations were added to the SSTs used for the atmospheric integrations and so the long-term mean of the heat fluxes are strongly constrained.

In all, 32 runs were conducted, each 15 years long, corresponding to the 1979-1993 period of the SST data used to force the atmospheric model. Five-day averages (pentads) of the model fields were archived. These were subsequently interpolated to the 11 depth levels, coincident with the depths of the TAO observations. All the covariance estimates have been made using these fields. Selecting at random a pentad from a 15-year period, a computation of the eofs of the matrix $\mathbf{P}$ was carried out using the ensemble of 32 ocean state realizations. The first eof explained only about 3% of the total error variance, and this result was similar for many one-time estimates of $\mathbf{P}$ attempted at other randomly selected dates. All eigenvalues of $\mathbf{AA}^T$ appear to be so close to each other as to be virtually indistinguishable. Apparently, this ensemble was not sufficient to reliably define the subspace containing the leading directions of the model error variability. A possible reason for this result is the small size of the ensemble, not adequate to resolve the dominant modes of variability of such a complex system. Thus, the question arose: how to enlarge the ensemble given the accumulated model output? A natural solution would be to include in the computation fields from the same model run, but selected in such a way as to prevent contamination of the internal model error variability by the temporal variability, such as lag correlation or interannual variations.

Thus, a matrix of ensemble members, $\mathbf{A}$, was formed by selecting at random five years from the 15 year period, then choosing a pentad from each year corresponding to the same date, say, the first of January. Such choice ensured that the states to be sufficiently separated in time to be considered independent. This allowed for collection of an ensemble of 160

members. This limit was set by practical considerations. The mean was removed separately for each of the 5 years to remove the influence of interannual variability. The eofs of the matrix **P** were then computed. The properties of the error covariance matrix constructed in such a way are discussed below.

*2)   Compact support*

A persistent problem associated with empirical forecast error covariance estimation is the appearance of unphysical large lag correlations that are an artifact of the limited ensemble size (e.g., Houtekamer and Mitchell 1998, fig. 6). We use an ensemble size of 160, yet the potential numbers of degrees of freedom are $O(10^6)$. To alleviate this problem, the multivariate anisotropic inhomogeneous matrix was modified by a matrix specified by a covariance function that vanishes at large distances; i.e., a Hadamard product of the two matrices was employed, as discussed by Houtekamer and Mitchell (2001). Keppenne and Rienecker (2002) implemented the compact support for the Ensemble Kalman Filter developed by the NASA Seasonal-to-Interannual Prediction Project (NSIPP) for parallel computing architectures, and that implementation is used in the present study. The functional form follows the work by Gaspari and Cohn (1999) who provided a methodology for constructing compactly supported multi-dimensional covariance functions. The characteristic scales of this function were selected in such a way that most of the local features of the empirically estimated error covariance structure are preserved, but at large spatial lags the covariance vanishes: 30°, 8° and 100 m in the zonal, meridional and vertical directions respectively.

To visualize the covariance structure, an artificial example is considered with a single observation different from a background field by one non-dimensional unit. The resulting correction reflects the model error correlation structure - it corresponds to a section of a single row of the **P** matrix. This is also termed the marginal gain since it measures the impacts of processing a single perfect measurement without reference to other data that might be assimilated. The correlation between temperature observations at several locations

across the equatorial Pacific ocean (156°E, 180°W, 155°W and 125°W) at depths roughly corresponding to the position of thermocline, estimated by the 20°C isotherm depth, and the temperature elsewhere in the Pacific reveals that with compact support the long range correlation is eliminated, but the local structure is intact (figure 3).

*3)  Multivariate error covariance patterns*

The following discussion of the multivariate error covariance model will focus on the thermocline region in the equatorial Pacific Ocean.  The shapes of the correlation structure associated with a single point differ between the eastern and western regions (figure 3, top 4 panels).  The zonal scale tends to be shorter in the western and central and longer in the eastern part of the basin. Meridional decay scales are similar along the equator, but the vertical correlation (figure 3, middle 4 panels) varies: shorter and symmetrical in the western part, slightly skewed in the central part and symmetrical but more elongated in the eastern part of the equatorial Pacific basin. Zonal sections (figure 3, bottom 4 panels) illustrate the anisotropy associated with the tilt of the thermocline.  This example alone demonstrates that the uniform temperature error covariance structure is so complex that a homogeneous error correlation structure is not quite applicable.

Although to date there have been very few salinity observations, this is changing with the Argo program (http://argo.jcommops.org, and Wilson, 2000).  Hence, it is of interest to explore corrections associated with salinity observations (figure 4). The decorrelation scales in the western basin are noticeably longer than in the middle and eastern basin, 8 to 10 degrees in zonal and 4 to 6 degrees in meridional direction in the west and 2-4 degrees in zonal and 1-2 degrees in meridional direction in the east. The scales are notably shorter than those for temperature (figure 3) except for the meridional scales in the west.

In a similar fashion one can analyze the temperature-salinity, temperature-velocity and other cross-variable relationships, i.e. the effect of a single unit observation on various fields - components of the ocean state vector. Corrections in S and U fields associated with a T

14

observation and corrections in T and U associated with an S observation are displayed for a single location, 155°W at equator (figure 5).

Examples of the temperature-salinity covariance (figure 5) reveal and reflect the complex and irregular nature of the temperature-salinity relationship. The change in salinity associated with a temperature increment is not necessarily density-compensating. Equatorial temperature and salinity south of the equator in the western region are anticorrelated, while temperature at the equator and salinity immediately to the north are correlated at 150 meters in the western and central Pacific. The scales of influence are short compared with the temperature-temperature relationship. The anticorrelation is consistent with the mean thermohaline (T-S) structure, with fresh water overlying a saline core. In the east, the correlation between T and S is primarily vertical; horizontal scales are very short, on the order of 2-4 degrees. The positive correlations on the equator, as seen on the meridional sections of the central basin, are higher towards the northern hemisphere. The negative correlations to the south are consistent with higher temperatures straddling the cold tongue with more saline water south of the equator and fresher water north. Thus the covariances are consistent with vertical and meridional variations.

The relationship between temperature and velocity in the western Pacific reflects temperature changes associated with upstream advection/convergence effects. At 156°E and at the dateline (not shown), the higher temperatures are associated with a weaker equatorial undercurrent in a broad region to the west. At 155°W, the effects are more local and wavelike with increased temperature associated with a stronger equatorial undercurrent. At 125°W (not shown) the scales are shorter and also wavelike, with changes in temperature apparently associated with instability waves.

It is possible to infer from the multivariate analysis the effect a single salinity observation would have on temperature and zonal velocity fields at various locations across the equatorial Pacific ocean. The high level of positive correlation between salinity and temperature field in the central and to a lesser degree in the eastern Pacific indicates that the correction of

the salinity field may have a significant impact on the temperature. The S-U relationship is weak in the western part of the basin and the correlation patterns are wavelike in the east, strongly pronounced in the north-south direction.

## 4.  Robustness of the model error covariance estimate

In this section, the sensitivity of the covariance structure to the choice made in populating the ensemble, i.e., to seasonal or interannual variations in the atmospheric forcing, is explored to evaluate the robustness of the covariance estimates. The robustness is tested by randomly sampling the full suite of integrations. Five years out of 15 (the length of the run) were picked at random, then the same date (e.g., January 1-5 pentad) was taken for each year. As before, the mean across the ensemble was removed for each year. The procedure was repeated ten times allowing us to obtain ten realizations of the covariance matrix $\mathbf{P}$. The pentads were chosen so that realizations from the same season and from different seasons could be compared. From visual assessment of figures similar to figures 3-5, the correlation structures represented by the different estimates of $\mathbf{P}$ were very similar.

One comparison of the robustness of covariance estimates is pointwise covariance sections (figure 6) at the same locations as simulated temperature observations as in figures 3-4. The tight distribution of the decorrelation curves from the 10 different $\mathbf{P}$ realizations (thin lines) indicates good reproducibility of the covariance structure. No significant interannual variability is apparent within this collection of $\mathbf{P}$ matrices. The over-plotted Gaussian curves show that the decorrelation scales vary at the four locations across the equatorial basin and can hardly be fitted by a single parameter (scale estimate) in a functional covariance model. In the UOI covariance model used for comparison below, the temperature decorrelation scales chosen are consistent with the scales of the empirical error covariance model in the western and central equatorial Pacific.

The difference among the Monte Carlo estimates of $\mathbf{P}$ can also be quantified in terms

16

of the dominant error subspaces spanned by each of the ensemble sets. These subspaces are best described by the orthonormal bases of empirical orthogonal functions (eofs). The use of eofs allows a spatial filtering of the covariance structures by inclusion of only those eofs that are non-noise-like, thus defining the dominant error subspace. This procedure also eliminates problems associated with different levels of variance even though the spatial structures (covariances) are similar.

Consider the projection of an ensemble of ocean state anomalies onto a given set of eofs. An anomalous ocean state vector $\mathbf{a}$ can be expressed in terms of the eof basis $\{\alpha\}$ as

$$\mathbf{a} = \Sigma_i a_i \alpha_i + \delta^\alpha. \tag{11}$$

The set of eofs $\{\alpha\}$ spans the subspace $\mathcal{S}_\alpha$ of the model error space $\mathcal{S}$ and $\delta^\alpha$ is the residual lying in the complement of $\mathcal{S}_\alpha$, i.e., subspace $\mathcal{S}_\alpha^c$, not spanned by $\{\alpha\}$. $\mathcal{S}_\alpha^c$ may or may not contain significant model error covariability information. To assess the information content not included in $\mathcal{S}_\alpha$ we examine covariability through the eofs of $\delta^\alpha$. If the eofs of $\delta^\alpha$ are noise-like, this would indicate that the eofs $\{\alpha\}$ captured the significant information regarding the model error covariance contained in $\mathbf{a}$. This calculation was repeated for several instances of $\{\alpha\}$ and $\mathcal{S} = \{\mathbf{a}\}$ to assess the invariability of $\mathcal{S}_\alpha$.

The spectra of various ensembles of $\delta^\alpha \subset \mathcal{S}_\alpha^c = \mathcal{S} \backslash \mathcal{S}_\alpha$ are shown in figure 7, where $\{\alpha\}$ are calculated from January pentads and $\{\mathbf{a}\}$ are pentads from July. In every case, the eigenvalues of $\{\alpha\}$ and $\{\delta\}$ are normalized by the variance of the corresponding ensemble $\{\mathbf{a}\}$. The eigencurves of $\{\delta\}$ are almost flat, characteristic of white noise, and are on order of magnitude less than the dominant eigenvalues of $\alpha$. Thus the error subspace generated from this Monte Carlo simulation appears to be robust.

## 5. Assimilation experiments

The effectiveness of the empirical multivariate forecast error covariance estimate is assessed by assimilating the temperature observations from the TAO moorings. The evaluation

uses a set of independent (i.e., not assimilated) temperature, salinity and zonal velocity observations from the TAO servicing cruises. The temperature and salinity data are based on Conductivity-Temperature-Depth (CTD) profiles and the velocity data from the Acoustic Doppler Profiler (ADCP). The comparison uses a gridded analysis of these data, as described by Johnson et al. (2000).

The assimilation experimental setup is as follows. The model was spun-up for 10 years with climatological forcing and then integrated with time dependent forcing for 1988-1998 in all the experiments. The assimilation began in July 1996. The initial conditions and the forcing were identical in all assimilation experiments. In addition to the data assimilation runs, a forced model integration without assimilation (referred to as the control) serves as a baseline for assessing the assimilation performance. The assimilation run with a simple univariate covariance model is denoted UOI. The run with the empirical multivariate forecast error covariance model is termed MvOI.

In every assimilation experiment, the daily-averaged subsurface temperature data from the TAO moorings were assimilated once a day. To alleviate the effects of the large shock on the model resulting from the intermittent assimilation of imperfectly balanced increments, the incremental update technique was used (Bloom et al. 1996). In this implementation, the assimilation increment is added gradually to the forecast fields at each time step.

The simulation (i.e., the control, with no assimilation) and two assimilation tests are cross-validated against the independent temperature, salinity and zonal velocity sections from Johnson et al. (2002). All of the available observed profiles are used and the statistics are separated corresponding to four regions: Niño 4 (160°E-150°W) and Niño 3 (150°W-90°W), further divided into two halves, south and north of the equator (0°-5°N and 5°S-0°). To put the amplitude of the RMSD in perspective, the mean monthly standard deviation (std) of the model is plotted as well. It is calculated using daily values at the same pre-defined depth levels on which the analyses are performed. The standard deviation represents the level of the internal variability in the model for the submonthly temporal scales which

could in part be responsible for the errors in the monthly averaged profiles assessed against single synoptic ship observations. In general, the RMSD of the control quantities and the data is about twice as large as the model standard deviation. The MvOI experiment shows comparable skill in temperature as the UOI with the greatest reduction in RMSD in the thermocline in the Niño 3 region south of the equator (figures 8 and 9). Below 400 meters neither of the assimilation schemes shows smaller RMSD than the control run due to the fact that data for assimilation are only available above 500 meters and at this level the observations are sparse. The transition between the upper part of the water column where the temperature profile is corrected by the assimilation to the abyss where the data are absent may cause disruptions in the internal dynamic balances. While the model is attempting to reinstate them using available mixing tools, it is not able to fully preserve the temperature structure below the transition region, which is reflected in the larger RMSD (top panels on figures 8 and 9). Apparently we should have calculated error covariance deeper to take care of this situation. [The problem has been corrected in the global implementation.] The MvOI is able, however, to preserve the salinity structure south of the equator and in the Niño 3 region north of the equator. To a lesser extent the MvOI current structure is also improved compared with the UOI, especially south of the equator.

The UOI assimilation improves upon the control case in the representation of temperature, yet the investigation of other model fields, such as salinity, reveals potential problems in a long-term integration. To illustrate this, consider time series of the equatorial salinity, averaged between 2°S and 2°N at the thermocline depth compared to the observed salinity (figure 10). In the UOI experiment, within 3-4 months the salinity structure deteriorates significantly. Poor performance of UOI is due to the fact that correcting the temperature field alone introduces artificial and potentially unstable water mass anomalies whose propagation and eventual strengthening destroys model dynamical balances. A method to alleviate this problem, proposed by Troccoli and Haines (1999) relies on the model-derived water mass properties to correct the model salinity commensurate with the temperature corrections made

19

by assimilating temperature observations. The salinity increments are calculated according to the temperature analysis by preserving the model's local T-S relationship. While the proposed method shows improvement in temperature and salinity analyses when tested with the Poseidon ocean model (Troccoli et al. 2003), it has the limitations that the scheme is designed solely for temperature observations and relies on the model maintaining a consistently good T-S relationship.

To test how well the assimilation schemes preserve the water mass properties, we consider, in a manner similar to Troccoli et al. (2003), the T-S relationships in the same subregions as used above. T-S pairs at each observation are compared with model values interpolated to the same locations using a T-S grid of granularity 0.25°C by 0.1 (figures 12 and 11). At least 5 T-S pairs must be found for a colored circle to be plotted to make sure that the features in the figures are robust. South and north of the equator in both Niño 3 and Niño 4 regions the model without assimilation (top panels) shows good representation of T-S except in the area of warmest water (cyan circles near the top of the plot) and somewhat in the representation of the dense cool saline water (few cyan circles below the main body of red color). The first deficiency is successfully corrected by the MvOI and to a lesser degree by the UOI. Some observed surface warm saline waters in the Niño 3 region north of the equator are not included in any of the model analyses, probably due to errors in surface forcing that the assimilation is not able to rectify. The problem of the lack of dense saline water in the model is slightly overcorrected by MvOI: all cyan circles change to red and some black circles appear in both regions north and south of the equator. The UOI scheme shows gross over-production of this type of water south of the equator and to a lesser degree in the north and it misses the more saline side of the distribution from $\sigma_\theta$ of 22 to 26 kg m$^{-3}$, north of the equator as well as in the south. Thus, significant problems are apparent in the UOI scheme, while MvOI is able to improve upon the control over almost the entire range of the T-S diagram.

Meridional cross-sections of the temperature, salinity and zonal velocity (figures 13, 14

and 15) are compared to a selection of sections prepared and presented in Johnson et al. (2002). The sections are chosen so that approximately simultaneous sections across the Pacific basin can be shown after a long period of integration (about 2 years). These sections are included in the RMSD statistics of figures 8 and 9. The temperature in the UOI experiment is an improvement over the control, while the salinity structure in the UOI has little resemblance to data. The model by itself is capable of producing good salinity and current fields. The UOI salinity cross sections display no penetration of the saline waters from the south across the equator and erroneous deep extension of high salinity around 2°S in the central and eastern basin. The MvOI salinity cross sections are more similar to the observations, although the salinity near the surface at 155°W north of the equator is somewhat low. The MvOI zonal current is the closest to the observed in the western and eastern Pacific with a better representation of the deeper subsurface maxima and a surfacing of the undercurrent at 165°E. The UOI currents reach too deep. At the dateline the current structure in MvOI is exaggerated compared to observed but the secondary subsurface maximum at about 4°N (the northern subsurface countercurrent) is captured in the assimilation. UOI currents are again too weak, particularly at the equator and reach too deep south of the equator. It is apparent from these figures that the MvOI corrects the current structure on and close to the equator better than the statistics of figures 8 and 9 might suggest.

## 6. Conclusion

Two conceptually different forecast error covariance models were considered in the context of the optimal interpolation data assimilation. One is the univariate model of the temperature error which uses a Gaussian spatial covariance function with different scales in zonal, meridional and vertical directions. The second is the multivariate error covariance matrix estimated in the dominant error subspace of empirical orthogonal functions (eofs)

generated from Monte Carlo simulations. The latter provides an empirical estimate of the covariability of the errors in temperature, salinity and current fields and spatial structure consistent with the governing dynamics. Thus during an assimilation cycle not only the temperature field, but the entire ocean state vector can be updated.

The univariate assimilation scheme brought the temperature field close to observations, yet the structure of the unobserved fields (salinity and currents) deteriorated quickly, precluding long-term integration. Most of the problems with the univariate OI run (no salty tongue in the south and deep penetration of salinity in the south, currents that are too deep) are due to neglect of the correlation between temperature and salinity when assimilating temperature alone which tends to cause spurious convective overturning. The multivariate scheme more successfully corrects the salinity and currents as verified by independent observations.

The empirical error covariance model presented in this study is an initial estimate of the forecast error covariance, and is used throughout the assimilation under the assumption that the forecast error statistics do not change significantly in time or after prior assimilation. The robustness of such an estimate was investigated and it was found that it does not exhibit significant seasonal or interannual variability, although there are not enough simulation years to distinguish among statistics during El Niño, La Niña and normal years.

The empirical multivariate forecast error covariance model provides important information regarding the error statistics of all the model fields, prognostic or diagnostic. This gives a natural way to include into the state estimation process observations of different types, for example, the sea surface height, which is often a model diagnostic.

Further developments are underway in implementing the MvOI method for the global ocean model configuration, particularly improving the ensemble statistics by including synoptic perturbations to the forcing fields, perturbations to the model parameters and initial conditions. It is more natural, taking into account the Poseidon ocean model formulation, to consider the covariances of the model variables within the quasi-isopycnal layers. Investigations are also underway to make the MvOI scheme more efficient in a reduced space by

including only a limited number of leading eofs.

## 7.  Acknowledgments

# REFERENCES

Atlas R., S. C. Bloom, R. N. Hoffman, J. V. Ardizzone, and G. Brin, 1991: Space-based surface wind vectors to aid understanding of air-sea interactions. *Eos, Transactions, Amer. Geophys. U.*, **72**, 201,204,205 and 208.

Bloom S. C., L. L. Takacs, A. M. D. Silva, and D. Ledvina, 1996: Data assimilation using incremental amalysis updates. *Mon. Wea. Rev.*, **124**, 1256–1271.

Borovikov A., M. M. Rienecker, and P. S. Schopf, 2001: Surface heat balance in the equatorial Pacific Ocean: Climatology and the warming event of 1994-95. *Journal of Climate*, **14(12)**, 2624–2641.

Buehner M., 2004: Ensemble-derived stationary and flow-dependent background error co-variances: Evaluation in a quasi-operational NWP setting. *Q. J. R. Meteorol. Soc.*, **128**, 1–29.

Burgers G., and P. J. van Leeuwen, 1998: Analysis scheme in the Ensemble Kalman Filter. *Mon. Wea. Rev.*, **126**, 1719–1724.

Cane M. A., A. Kaplan, R. N. Miller, B. Tang, E. C. Hackert, and A. J. Busalacchi, 1996: Mapping tropical Pacific sea level: Data assimilation via a reduced state Kalman filter. *J. of Geophys. Res.*, **101(C10)**, 22599–22617.

Carton J. A. and E. C. Hackert, 1990: Data assimilation applied to the temperature and circulation in the tropical Atlantic, 1983-1984. *J. of Phys. Oceanogr.*, **20(8)**, 1150–1164.

Cohn S. E., 1997: An introduction to estimation theory. *Journal of the Meteorological Society of Japan*, **75**, 257–288. Special issue dedicated to "Data Assimilation in Meteorology and Oceanography: Theory and Practice".

Daley R., 1991: *Atmospheric Data Analysis.* Cambridge University Press, 457 pp.

Derber J. A., D. F. Parrish, and S. J. Lord, 1991: The new global operational system at the National Meteorological Center. *Wea. and Forecast*, **6**, 538–547.

Evensen G., 1994: Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics. *J. Geophys. Res.*, **C99**, 10143–10162.

Fisher M., and E. Andersson, 2001: Developments in 4D-Var and Kalman filtering. Technical report, ECMWF.

Freitag H. P., Y. Feng, L. J. Mangum, M. P. McPhaden, J. Neander, and L. D. Stratton, 1994: Calibration procedures and instrumental accuracy estimates of TAO temperature, relative humidity and radiation measurements. ERL PMEL-104, NOAA.

Fukumori I., J. Benveniste, C. Wunch, and D. B. Haidvogel, 1993: Assimilation of sea surface topography into an ocean circulation model using a steady-state smoother. *J. of Phys. Oceanography*, **23**, 1831–1855.

Gaspari G., and S. E. Cohn, 1999: Contruction of correlatoin functions in two and three dimensions. *Quart. J. Roy. Meteorol. Soc.*, **125**, 723–757.

Ghil M., and P. Malanotte-Rizzoli, 1991: Data assimilation in meteorology and oceanography. *Advances in Geophysics*, **33**, 141–266.

Harrison D. E., B. S. Giese, and E. S. Sarachik, 1990: Mechanisms of SST change in the equatorial waveguide during the 1982-83 ENSO. *J. Clim.*, **3**, 173–188.

Houtekamer P. L., L. Lefaivre, J. Derome, H. Ritchie, and H. L. Mitchell, 1996: A system simulation approach to ensemble prediction. *Mon. Wea. Rev.*, **124**, 1225–1242.

———, and H. Mitchell, 1998: Data assimilation using an ensemble Kalman filter technique. *Mon. Wea. Rev.*, **126**, 796–811.

‗‗‗‗‗‗, and H. Mitchell, 2001: A sequential ensemble Kalman filter for atmospheric data assimilation. *Mon. Wea. Rev.*, **129**, 123–137.

Ji M., A. Leetmaa, and J. Derber, 1995: An ocean analysis system for seasonal to interannual climate studies. *Mon. Wea. Rev.*, **123**, 460–481.

Johnson G. C., M. J. McPhaden, G. D. Rowe, and K. E. McTaggart, 2000: Upper equatorial Pacific Ocean current and salinity variability during the 1996-1998 El Niño-La Niña cycle. *J. of Geophys. Res.*, **105**, 1037–1053.

Kalman R., 1960: A new approach to linear filtering and prediction problems. *J. Basic Eng.*, **D82**, 35–45.

Kalnay E., and coauthors, 1996: The NCEP/NCAR 40-year reanalysis project. *Bulletin of the Amer. Meteorol. Soc.*, **77**, 437–471.

Kaplan A., Y. Kushnir, M. A. Cane, and M. B. Blumenthal, 1997: Reduced space optimal analysis for historical data sets: 136 years of Atlantic sea surface temperature. *J. of Geophys. Res.*, **102**, 27835–27860.

Keppenne C. L., and M. M. Rienecker, 2002: Initial testing of a parallel ensemble Kalman filter with the Poseidon isopycnal ocean general circulation model. *Mon. Wea. Rev.*, **130**, 2951–2965.

‗‗‗‗‗‗, and M. M. Rienecker, 2003: Assimilation of temperature into an isopycnal general circulation using a parallel Ensemble Kalman Filter. *J. Mar. Sys.*, **40-41**, 363–380.

Kraus E. B., and J. S. Turner, 1967: A one-dimensional model of the seasonal thermocline: II. the general theory and its consequences. *Tellus*, **19**, 98–109.

Large W. G., and S. Pond, 1982: Open ocean momentum flux measurements in moderate to strong winds. *J. Phys. Oceanogr.*, **11**, 324–336.

Lorenc A. C., 1986: Analysis methods for numerical weather prediction. *Quart. J. R. Met. Soc.*, **112**, 1177–1194.

McPhaden M. J., A. J. Busalacchi, R. Cheney, J.-R. Donguy, K. S. Gage, D. Halperin, M. Ji, P. Julian, G. Meyers, G. T. Mitchum, P. P. Niiler, J. Picaut, R. W. Reynolds, N. Smith, and K. Takeuchi, 1998: The tropical ocean global atmosphere observing system: A decade of progress. *J. of Geophys. Res.*, **103**, 14169–14240.

Oke P. R., J. S. Allen, R. N. Miller, G. D. Egbert, and P. M. Kosro, 2002: Assimilation of surface velocity data into a primitive equation coastal ocean model. *J. of Gephys . Res.*, **107(C9)**, 5–1–5–25.

Pacanowski R. and S. G. H. Philander, 1981: Parametrization of vertical mixing in numerical models of tropical oceans. *J. Phys. Oceanogr.*, **11**, 1443–1451.

Philander S. G , 1990: *El Niño, La Niña, and the Southern Oscillation.* Academic Press, 239 pp.

Preisendorfer R. W., 1988: *Principal component analysis in meteorology and oceanography.* Elsevier, 425 pp.

Rienecker M. M., and R. N. Miller, 1991: Ocean data assimilation using optimal interpolation with a quasi-geostrophic model. *Journal of Geophysical Research*, **96(C8)**, 15093–15103.

Rosati A., R. Gudgel, and K. Miyakoda, 1996: Global ocean data assimilation system. *Modern Approaches to Data Assimilation in Ocean Modeling*, Elsevier.

——, K. Miyakoda, and R. Gudgel, 1997: The impact of ocean initial conditions on ENSO forecasting with a coupled model. *Mon. Wea. Rev.*, **125**, 754–772.

Rossow W. B., and R. A. Schiffer, 1991: ISCCP cloud data products. *Bull. Am. Met. Soc.*, **72**, 2–20.

Schopf P. S., and A. Loughe, 1995: A reduced-gravity isopycnal ocean model - hindcasts of El Niño. *Mon. Wea. Rev.*, **123**, 2839–2863.

Seager R., M. B. Blumenthal, and Y. Kushnir, 1994: An advective atmospheric mixed layer model for ocean modeling purposes: Global simulation of surface heat fluxes. *J. Clim.*, **8**, 1951–1964.

Shapiro R., 1970: Smoothing, filtering and boundary effects. *Rev. of Geophys. Space Phys.*, **8**, 359–387.

Suarez M. J., and L. L. Takacs, 1995: Documentation of the Aries/GEOS dynamical core Version 2. Technical memorandum 104606, **5**, NASA. pp. 44.

Troccoli A., and K. Haines, 1999: Use of the temperature-salinity relation in a data assimilation context. *J. Atmos. Oceanic Technol.*, **16**, 2011–2025.

_____, M. A. Balmaseda, J. Segschneider, J. Vialard, D. L. T. Anderson, K. Heines, T. Stockdale, F. Vitart, and A. D. Fox, 2002: Salinity adjustments in the presence of temperature data assimilation. *Mon. Wea. Rev.*, **130**, 89–102.

_____, M. M. Rienecker, C. L. Keppenne, and G. C. Johnson, 2003: Temperature data assimilation with salinity corrections: Validation in the tropical Pacific Ocean, 1993-1998. Technical Report 104606, **24**, NASA GSFC. pp. 23.

Wilson S., 2000: Launching the Argo armada. *Oceanus*, **42**, 17–19.

Xie P. P., and P. Arkin, 1997: Global precipitation: a 17-year monthly analysis based on gauge observations, satellite estimates, and numerical model outputs. *Bull. Amer. Meteorol. Soc.*, **11**, 2539–2558.

Yang S., K.-M. Lau, and P. S. Schopf, 1999: Sensitivity of the tropical Pacific Ocean to precipitation-induced freshwater flux. *Climate Dynamics*, **15**, 737–750.

Zhang S., and J. L. Anderson, 2003: Impact of spatially and temporally varying estimates of error covariance on assimilation in a simple atmospheric model. *Tellus*, **55A**, 126–147.

# LIST OF FIGURES
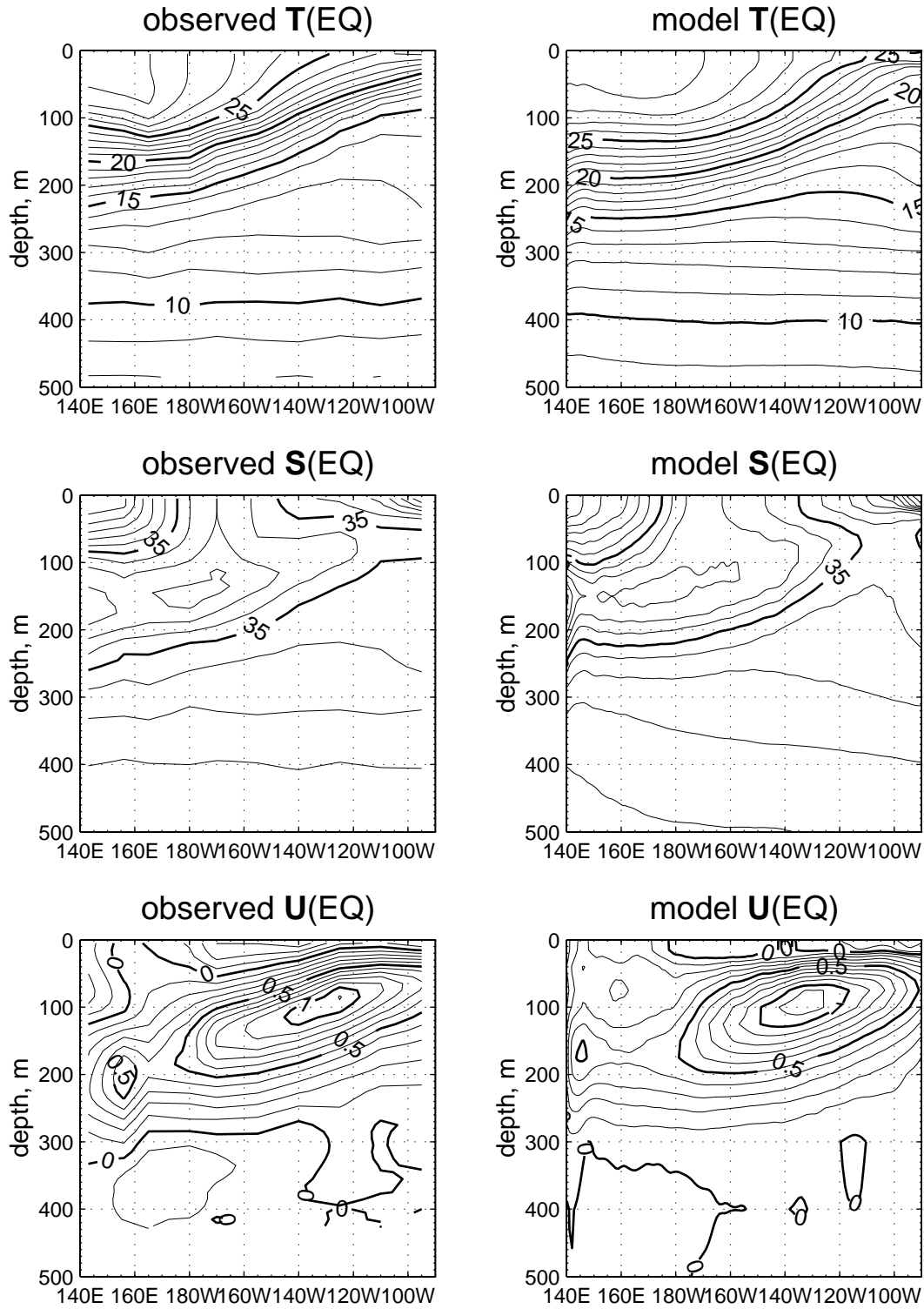
FIG. 1. Equatorial cross-section of the Poseidon model means (1988-1997) of temperature, salinity and zonal velocity (right panels) and corresponding data-based estimates (left panels) from Johnson et al.(2002).
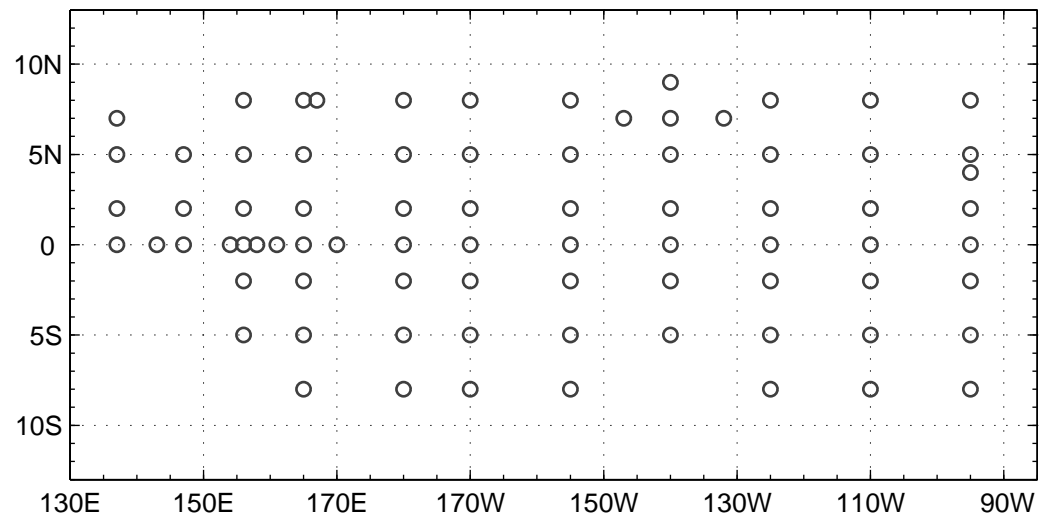
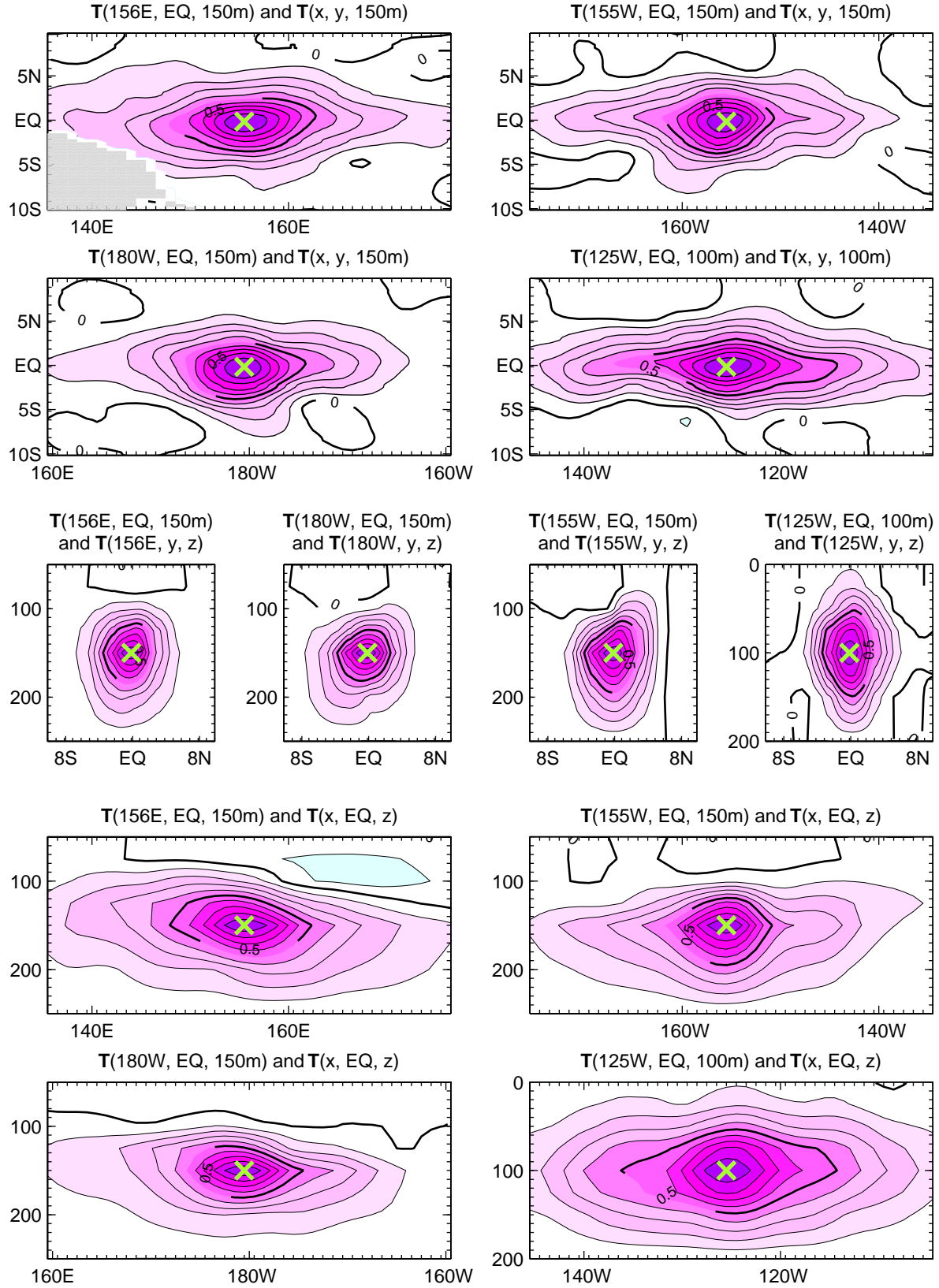FIG. 2. Map of the TAO array, consisting of approximately 70 moored ocean buoys in the Tropical Pacific Ocean.

FIG. 3. Examples of correlation structure derived from a 160 member ensemble. The compact support is applied as described in the text. Contour interval is 0.1. Crosses mark the position of the simulated observation.
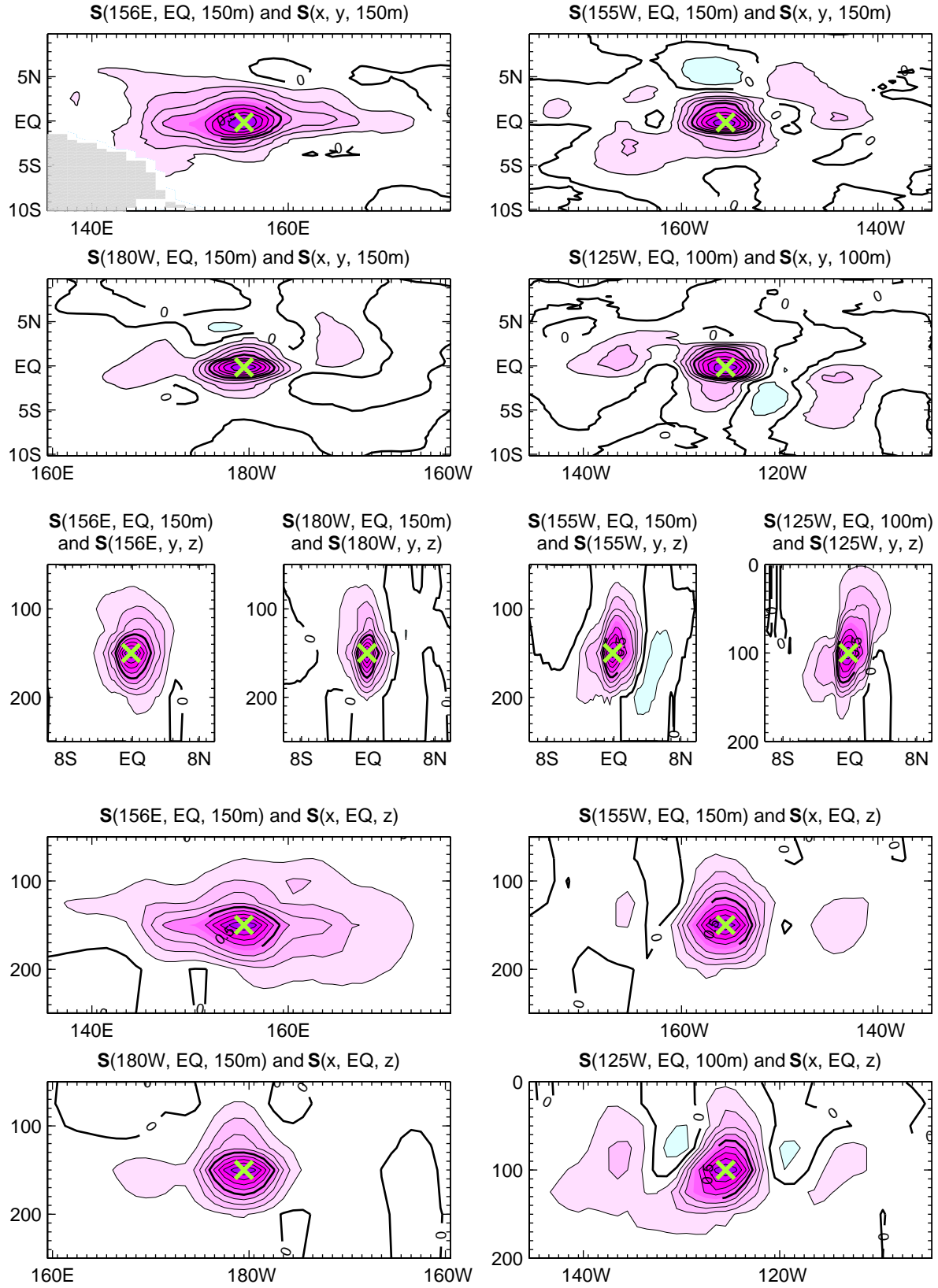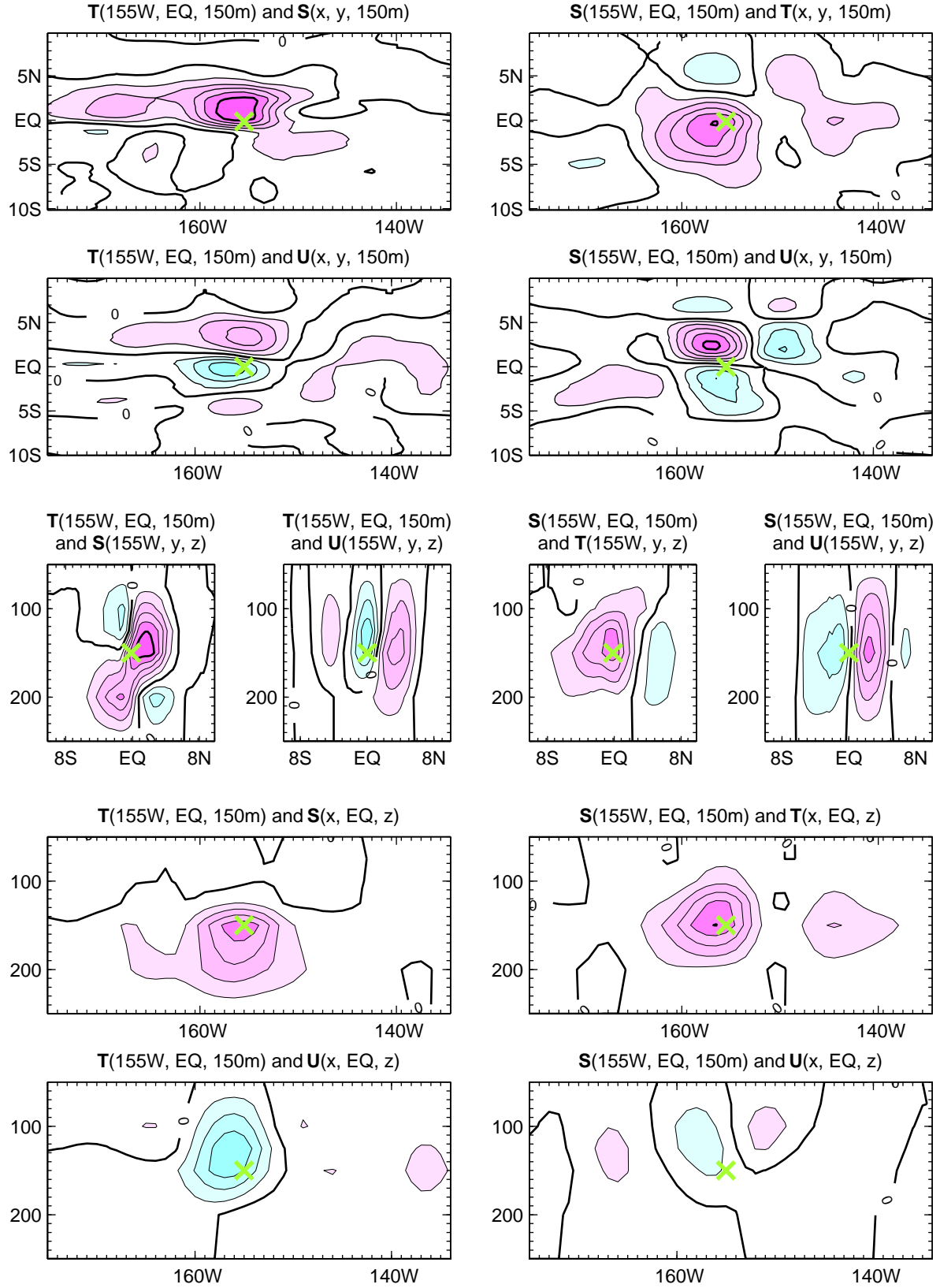
FIG. 4. Same as in figure 3 but for salinity.

FIG. 5. Examples of correlation structure derived from a 160 member ensemble. The compact support is applied as described in the text. Various combinations of observed and updated variables are presented. Contour interval is 0.1. Crosses mark the position of the simulated observation.

FIG. 6. One-dimensional decorrelation curves (zonal, meridional and vertical directions) corresponding to the simulated observation at the specified locations. Each thin solid line produced by a different realization of the error covariance matrix. Dashed grey lines show the Gaussian functional error covariance model used in UOI.

FIG. 7. Eigenvalues for several realizations of the matrix $\mathbf{P}$ (marked $\alpha$) and the eigenvalues for ensembles of $\delta$'s - the residuals of the projections of an arbitrary collection of anomalous ocean states onto a basis of eofs.

FIG. 8. RMSD between the three model runs (UOI, MvOI and control) and the observations as a function of depth for the 35 transects. Statistics are grouped by Niño 4 (160°E-150°W) and Niño 3 (150°W-90°W) regions, and each area is further divided into two halves, south and north of the equator (0°-5°N) shown here). Temperature RMSD (a-b), salinity RMSD (c-d) and zonal velocity RMSD (e-f) are shown. Mean monthly standard deviations of the corresponding model fields for the same regions are shown by stars.

FIG. 9. As in figure 8, but for regions south of the equator (5°S-0°).

FIG. 10. Salinity time series for the control, UOI and MvOI integrations. CTD observations are shown where available. Values are averaged between 2°S-2°N at the specified longitudes.

FIG. 11. Temperature-Salinity diagram for UOI, MvOI and control experiments for Niño 4 and Niño 3 regions south of the equator. Black dot is plotted for values present only in the model, cyan - only in observations and points where the model and observations agree are shown in red.

FIG. 12. As in figure 11, but for regions south of the equator (5°S-0°).

FIG. 13. Meridional profiles of the model and observed temperature. Model fields are averaged over 1 month, whereas the observations are from individual quasi-synoptic CTD/ADCP sections (following Johnson et al. 2000).

FIG. 14. As for figure 13, but for salinity. Contour interval is 0.2.

FIG. 15. As in figure 13, but for zonal velocity. Contour interval is 0.2 ms$^{-1}$.

Figure 1: Equatorial cross-section of the Poseidon model means (1988-1997) of temperature, salinity and zonal velocity (right panels) and corresponding data-based estimates (left panels) from Johnson et al.(2002).

Figure 2: Map of the TAO array, consisting of approximately 70 moored ocean buoys in the Tropical Pacific Ocean.

Figure 3: Examples of correlation structure derived from a 160 member ensemble. The compact support is applied as described in the text. Contour interval is 0.1. Crosses mark the position of the simulated observation.

Figure 4: Same as in figure 3 but for salinity.

Figure 5: Examples of correlation structure derived from a 160 member ensemble. The compact support is applied as described in the text. Various combinations of observed and updated variables are presented. Contour interval is 0.1. Crosses mark the position of the simulated observation.

Figure 6: One-dimensional decorrelation curves (zonal, meridional and vertical directions) corresponding to simulated observation at the specified locations. Each thin solid line is produced by a different realization of the error covariance matrix. Dashed grey lines show the Gaussian functional error covariance model used in UOI.

Figure 7: Eigenvalues for several realizations of the matrix **P** (marked $\alpha$) and the eigenvalues for ensembles of $\delta$'s - the residuals of the projections of an arbitrary collection of anomalous ocean states onto a basis of eofs.

Figure 8: RMSD between the three model runs (UOI, MvOI and control) and the observations as a function of depth for the 35 transects. Statistics are grouped by Niño 4 (160°E-150°W) and Niño 3 (150°W-90°W) regions, and each area is further divided into two halves, south and north of the equator (0°-5°N) shown here. Temperature RMSD (a-b), salinity RMSD (c-d) and zonal velocity RMSD (e-f) are shown. Mean monthly standard deviations of the corresponding model fields for the same regions are shown by stars.

Figure 9: As in figure 8, but for regions south of the equator (5°S-0°).

Figure 10: Salinity time series for the control, UOI and MvOI integrations. Values are averaged between 2°S-2°N at the specified longitudes. CTD observations are shown by stars where available.
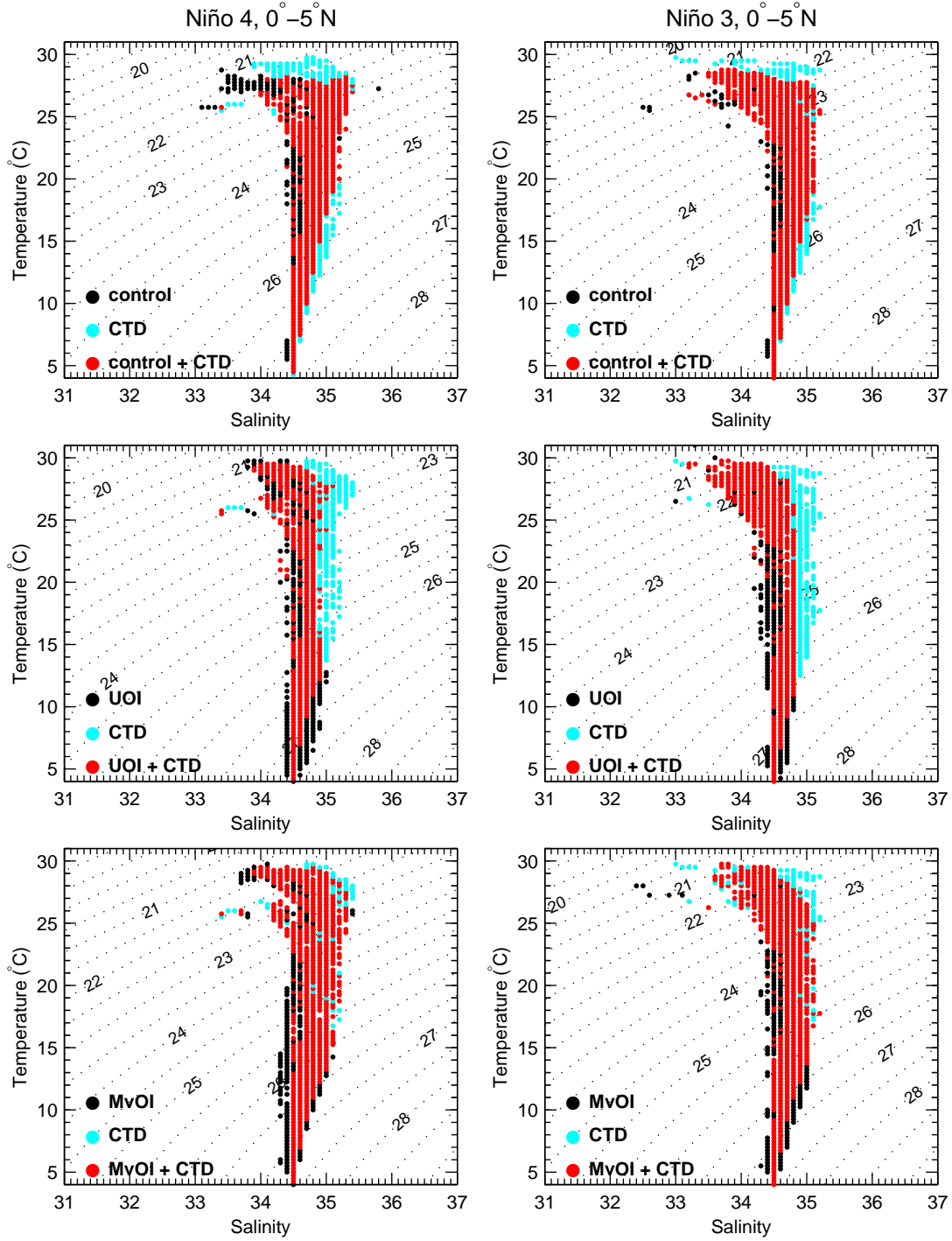
Figure 11: Temperature-Salinity diagram for UOI, MvOI and control experiments for Niño 4 and Niño 3 regions north of the equator. Black dots are plotted for values present only in the model, cyan - only in observations and points where the model and observations agree are shown in red.

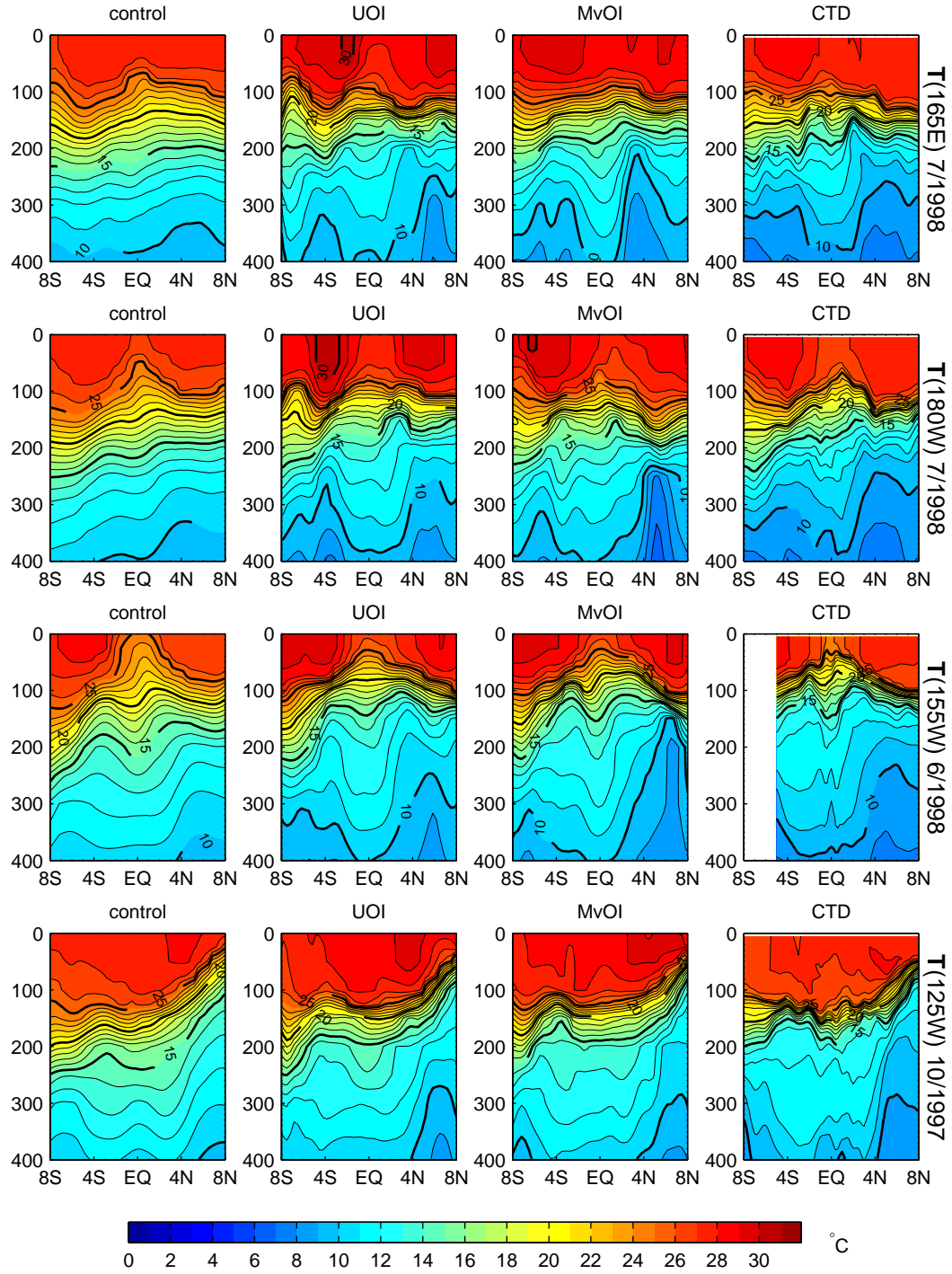Figure 12: As in figure 11, but for regions south of the equator (5°S-0°).

Figure 13: Meridional vertical sections of the model and observed temperature. Model fields are averaged over one month, whereas the observations are from individual quasi-synoptic CTD/ADCP sections (following Johnson et al. 2000). Contour interval is 1°C.
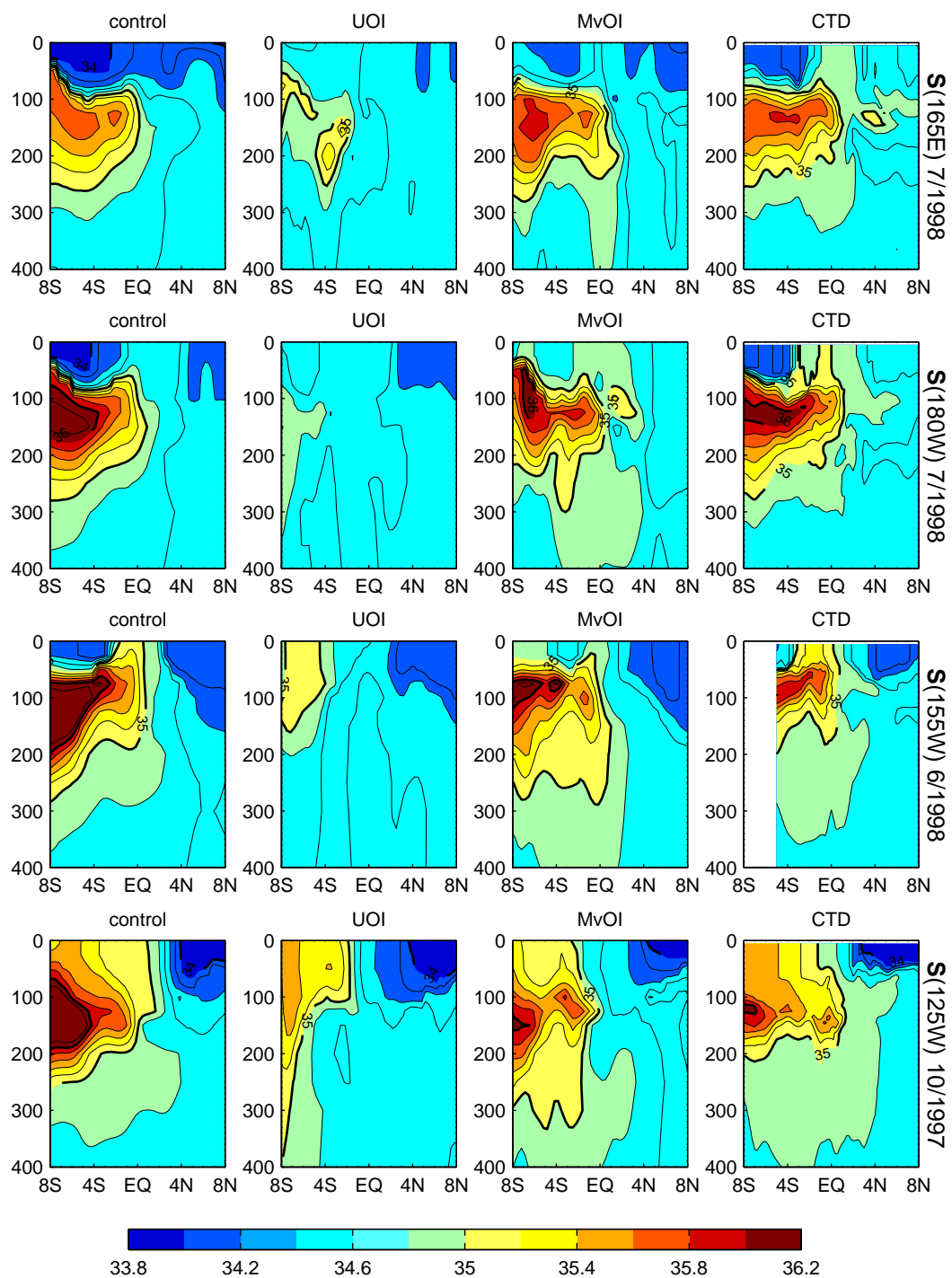
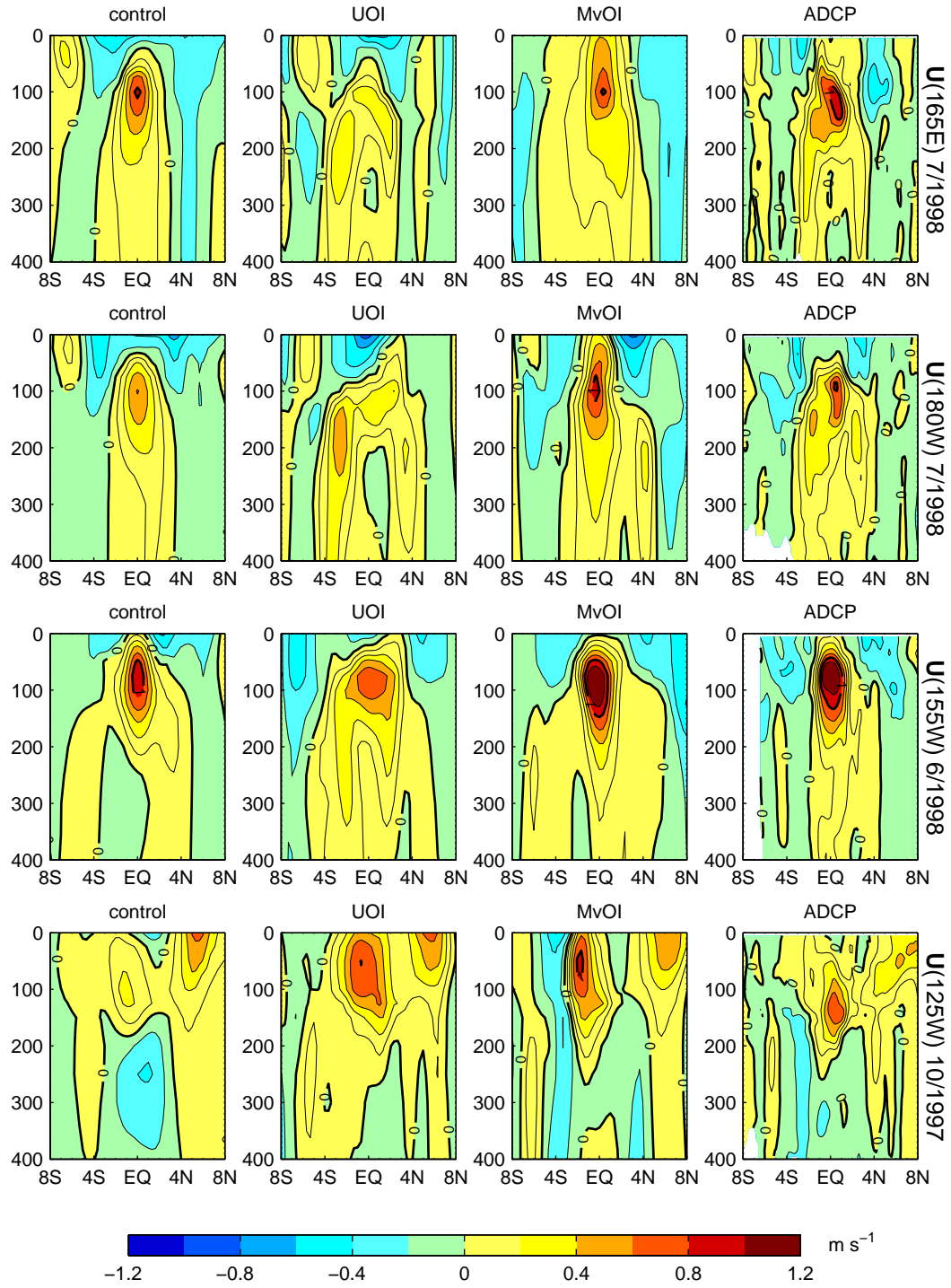Figure 14: As for figure 13, but for salinity. Contour interval is 0.2.

Figure 15: As in figure 13, but for zonal velocity. Contour interval is $0.2 \text{ ms}^{-1}$.